

## PATENT ABSTRACTS OF JAPAN

2

(11)Publication number : 09-185392

(43)Date of publication of application : 15.07.1997

(51)Int.Cl.

G10L 3/02

G10H 1/00

G10K 15/04

(21)Application number : 07-353508

(71)Applicant : VICTOR CO OF JAPAN LTD

(22)Date of filing : 28.12.1995

(72)Inventor : NIIHARA TOSHIKO  
MATSUMOTO MITSUO  
SUZUKI TAKUMA

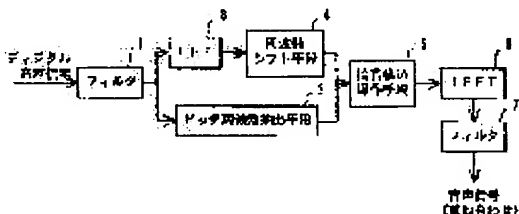
## (54) INTERVAL CONVERTING DEVICE

## (57)Abstract:

PROBLEM TO BE SOLVED: To convert the interval of an individual's voice without deterioration in sound quality so that features of the individual's voice are left.

SOLUTION: A digital input voice signal is cut by a filter 1 into frames of a specific time and a pitch frequency extracting means 2 extracts the pitch frequency of the voice signal outputted from this filter

1. The voice signal outputted from the filter 1 is supplied to an FFT(fast Fourier transforming means) circuit 3 as well and converted from a time-area signal from a frequency-area signal, whose entire frequency band is shifted by a frequency shift means 4 to a higher or lower frequency side. Then a harmonic structure operating means 5 increases or decreases the level of harmonic components of the pitch frequency of the voice signal having its entire frequency band shifted by the frequency shift means 4 from the pitch frequency extracted by the pitch frequency extracting means 2. Then an IFFT(inverse fast Fourier transforming means) circuit 6 converts the signal into a time-area signal, which is outputted.



## LEGAL STATUS

[Date of request for examination] 24.03.2000

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number] 3265962

[Date of registration] 11.01.2002

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's  
decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2000 Japan Patent Office

2

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平9-185392

(43) 公開日 平成9年(1997)7月15日

(51) Int.Cl. <sup>6</sup>	識別記号	庁内整理番号	F I	技術表示箇所
G 1 0 L 3/02			G 1 0 L 3/02	A
G 1 0 H 1/00			G 1 0 H 1/00	B
G 1 0 K 15/04	3 0 2		G 1 0 K 15/04	3 0 2 D

審査請求 未請求 請求項の数 3 F D (全 7 頁)

(21) 出願番号 特願平7-353508

(22) 出願日 平成7年(1995)12月28日

(71) 出願人 000004329

日本ビクター株式会社

神奈川県横浜市神奈川区守屋町3丁目12番地

(72) 発明者 新原 寿子

神奈川県横浜市神奈川区守屋町3丁目12番地 日本ビクター株式会社内

(72) 発明者 松本 光雄

神奈川県横浜市神奈川区守屋町3丁目12番地 日本ビクター株式会社内

(72) 発明者 鈴木 琢磨

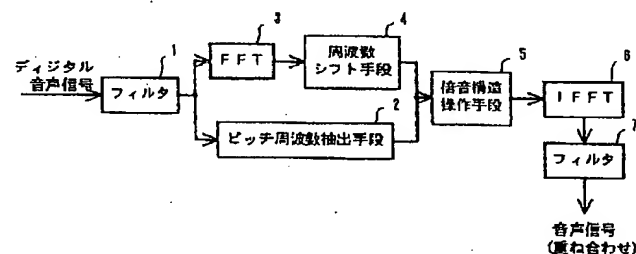
神奈川県横浜市神奈川区守屋町3丁目12番地 日本ビクター株式会社内

(54) 【発明の名称】 音程変換装置

(57) 【要約】

【課題】 音質劣化がなく、個人の声の特徴を残した音声の音程変換ができなかった。

【解決手段】 デジタル入力された音声信号はフィルタ1により所定時間のフレームに切り出され、このフィルタ1から出力される音声信号のピッチ周波数をピッチ周波数抽出手段2により抽出する。また、フィルタ1から出力される音声信号は、FFT回路3にも供給され、時間領域の信号から周波数領域の信号へ変換した後、周波数シフト手段4により全周波数帯域を高域側または低域側にシフトする。そして、倍音構造操作手段5は、ピッチ周波数抽出手段2により抽出されたピッチ周波数から、周波数シフト手段4により全周波数帯域をシフトされた音声信号のピッチ周波数の倍音成分のレベルを増減させる。その後、IFFT回路6により時間領域の信号に変換して出力する。



**【特許請求の範囲】**

**【請求項 1】** デジタル入力された音声信号を所定時間の時間窓で切り出す分割手段と、この分割手段から出力される音声信号の基本周波数を抽出するピッチ周波数抽出手段と、前記分割手段から出力される音声信号を時間領域の信号から周波数領域の信号へ変換するフーリエ変換手段と、このフーリエ変換手段より出力される音声信号の全周波数帯域を高域側または低域側にシフトする周波数シフト手段と、前記ピッチ周波数抽出手段により抽出されたピッチ周波数が供給され、前記周波数シフト手段により全周波数帯域をシフトされた音声信号の倍音の構造を操作する倍音構造操作手段と、この倍音構造操作手段より出力される音声信号を時間領域の信号に変換する逆フーリエ変換手段とを有することを特徴とする音程変換装置。

**【請求項 2】** 前記分割手段は、デジタル入力された音声信号を所定時間のフレームに切り出すと共に、このフレームの最初の部分の 10～35 msec のデータを正弦波の 1/4 周期分の時間窓で切り出し、このフレームの最後の部分の 10～35 msec のデータを余弦波の 1/4 周期分の時間窓で切り出すことを特徴とする請求項 1 記載の音程変換装置。

**【請求項 3】** 前記倍音構造操作手段は、前記全帯域シフト手段により高域側へシフトされた際には音声信号の倍音成分のレベルを減少させ、低域側へシフトされた際には音声信号の倍音成分のレベルを増加させることを特徴とする請求項 1 または請求項 2 記載の音程変換装置。

**【発明の詳細な説明】****【0001】**

**【発明の属する技術分野】** 本発明は、カラオケ装置や音響映像編集装置等に使用され、音声の音程（ピッチ周波数、基本周波数）を変換する音程変換装置に係り、特に、音質の劣化がなく、かつ個人の声の特徴を残したまま音声の音程を容易に変換することのできる音程変換装置に関するものである。

**【0002】**

**【従来の技術】** 従来より、カラオケ装置等では、歌う人の音域に合わせるために、演奏される伴奏の音程を自由に变化させて設定することができるキーコントロールと呼ばれる機能が付いていた。これは、伴奏として再生されるアナログ音声信号の再生速度を変化させることにより、音程を変化させていた。また、近年では、センタに曲のデータを蓄積しておき、このセンタに複数接続されている遠隔地の端末装置に必要な応じて曲のデータを送信して、端末装置で曲を再生する通信カラオケが開発されている。

**【0003】** この通信カラオケのセンタから端末装置に送信される曲のデータは、曲に合わせて歌詞を表示する

と共にその表示色を変更するための文字データと、曲の伴奏を再生するために端末装置のシンセサイザを動作させる MIDI 信号と、男性または女性の声による肉声バックコーラスを端末装置で再生するための圧縮された音声信号とで構成されている。そして、この通信カラオケの端末装置において、演奏される伴奏の音程を変える場合、MIDI 信号で再生されるシンセサイザの音程を、全体的に上げる（下げる）様に設定することにより、再生速度を変えずに音程を自由に変えて再生することができる。

**【0004】** ところが、肉声バックコーラスは、MIDI 信号でないため、音程に関連するデータを備えておらず、再生速度を変えない状態で、音質の劣化がなく、しかも個人の声の特徴を残したまま音声の音程を変換することは困難であった。また、近年の音響映像編集装置は、デジタル信号の状態編集作業を行うものも開発されてきているが、高品質を維持したまま音声の音程を変換させるのは困難であった。

**【0005】** これまでの音声の再生速度を一定に保ったまま音声の音程を変換する方法としては、主として二通りの方法が考えられている。一つは、音声波形を時間領域で操作する方法であり、例えばピッチ周波数を 2 倍に上げる場合、音声信号を所定時間毎に切り出して、この切り出し区間毎に 2 倍の速度でデータを読み出すようにしている。そしてこの場合、切り出した区間のデータからピッチ周波数（ピーク周波数のうち最も低い周波数）を求め、2 倍のピッチ周波数である波形を付け加えることで時間を変えずにピッチ周波数のみ 2 倍に上げることができる。さらに、この様な処理をした切り出し区間をスムーズに繋げることによって音程変換を実現することができるが、現実には、繋げ方によって音質を損ねたり、個人の声の特徴が維持されず不自然な音声になってしまうので、現在も各種改善方法が提案されている状態である。

**【0006】** もう一つは、フーリエ変換を用いて周波数領域で操作する方法である。音声信号を所定時間毎に切り出し、フーリエ変換によって周波数の振幅成分と周波数の位相成分とを抽出する。次に、全周波数帯域を所望のシフト量分だけ周波数シフト及び位相シフトし、逆フーリエ変換した後、切り出し区間を繋げていく方法である。しかし、この方法によっても不自然な音声となってしまう、うまく音程変換ができなかった。なお、フーリエ変換後、ピークスpekトル（ピッチ周波数）を検出し、このピークスpekトル付近の周波数信号のみをシフトする方法が当社より出願され、特開昭 59-204096 号公報に公開されている。

**【0007】**

**【発明が解決しようとする課題】** 特開昭 59-204096 号公報に記載されている、ピークスpekトルを示す周波数成分のみシフトを行なう方法は、ピークスpekトル

3

ルの倍音成分がそのまま残っているため、聴覚において元の音程が容易に想像されてしまい、倍音成分による元の音程とシフトした後の音程との2重の音程が聴こえてしまうという課題があった。

【0008】また、VTRやテープレコーダ等において、解説やナレーション等の音声を高速再生する際に、高くなってしまうピッチ周波数を元にもどして、聞き取りやすくするなど、カラオケのキーコントロール以外にも、音声のピッチ周波数を自由に変換したいという要求があった。そこで本発明は、従来に比べ簡単な回路構成で処理時間も比較的短く、しかも音質の劣化がなくて個人の声の特徴を維持したままの自然な音声音程変換を可能とする高品質な音程変換装置を提供することを目的とする。

【0009】

【課題を解決するための手段】本発明は、上記目的を達成するための手段として、デジタル入力された音声信号を所定時間の時間窓で切り出す分割手段と、この分割手段から出力される音声信号の基本周波数を抽出するピッチ周波数抽出手段と、前記分割手段から出力される音声信号を時間領域の信号から周波数領域の信号へ変換するフーリエ変換手段と、このフーリエ変換手段より出力される音声信号の全周波数帯域を高域側または低域側にシフトする周波数シフト手段と、前記ピッチ周波数抽出手段により抽出されたピッチ周波数が供給され、前記周波数シフト手段により全周波数帯域をシフトされた音声信号の倍音の構造を操作する倍音構造操作手段と、この倍音構造操作手段より出力される音声信号を時間領域の信号に変換する逆フーリエ変換手段とを有することを特徴とする音程変換装置を提供しようとするものである。

【0010】

【発明の実施の形態】以下、添付図面を参照して本発明の音程変換装置の一実施例を説明する。図1は本発明の音程変換装置の一実施例を示すブロック図であり、図2はその動作を示すフローチャート図である。そして、サンプリング周波数44.1kHzのデジタル音声信号が入力され、この音声信号を3半音高い方へピッチシフトする（音程を上げる）場合を例にして、以下に説明する。

【0011】まず、フレーム（処理区間）の番号（i）を初期化しておく（ステップ11）。そして、デジタル入力される音声信号がこのフレームよりも大きければ（ステップ12→Yes）、フィルタ（分割手段）1により4096サンプル毎のフレームに区切られて読み出され（ステップ13）、そのうち第0番～第999番のサンプル（最初の部分）は正弦波の窓関数で切り出され、第3096番～第4095番のサンプル（最後の部分）は余弦波の窓関数で切り出され、その他のサンプルは1の窓関数で切り出されて出力される（ステップ14）。なお、この正弦波及び余弦波の窓関数による時間窓での

4

切り出しは、後述する切り出し区間の重ね合わせの際に重ね合わせ部分の電力を一定にして各フレームをスムーズに繋げるために行うものである（図3参照）。

【0012】そして、このフィルタ1における正弦波および余弦波による時間窓での切り出しは、200～2000サンプル幅の任意サンプル幅の区間で種々実験したところ、音源によって多少の変化はあるが、ほとんどの音源で500～1500サンプル（約10～35ms）幅の間に最適な区間になることが判ったので、この実施例では1000サンプル（約23ms）幅で正弦波および余弦波による時間窓での切り出しを行っている。このフィルタ1により切り出された音声信号は、ピッチ周波数抽出手段2に供給されて、自己相関関数やケプストラム法等によりピッチ周波数（ピーク周波数のうち最も低い周波数（基本周波数）を示すサンプル）が抽出される（ステップ15）。また、フィルタ1より出力された音声信号は、FFT回路（フーリエ変換手段）3にも供給されてフーリエ変換を施され、時間領域の信号から周波数領域の信号へ変換される（ステップ16）。

【0013】このとき、時間領域に対応していた各サンプルは、各周波数に対応し、サンプル番号と周波数とが対応することになる。即ち、サンプリング周波数 $f_s$ の音声信号データをN個のサンプル毎に切り出して処理する場合、FFT回路3から出力される信号の周波数 $pHz$ を示すサンプル番号は第 $(p \times N / f_s)$ 番目となる。本実施例の場合、サンプリング周波数44.1kHzの音声信号データに対して4096サンプル毎に切り出しているため周波数 $pHz$ を示すサンプル番号は第 $(p \times 4096 / 44100)$ 番目となる（小数点以下切り捨て）。

【0014】そして、周波数シフト手段4により、実部と虚部とをピッチシフト量（3半音分）だけ移動させる（ステップ17）。ここで、1オクターブ（12半音）高い方へ移動させるということは、周波数を2倍にすることと同意であるので、h半音上げるには全体の周波数を $2^{h/12}$ 倍に上げれば良いことになる。ここでは、3半音高い方へずらすので、全体の周波数を $2^{3/12}$ 倍（約1.19倍）にすれば良い。その結果、第n番目のサンプルの値は第 $(1.19 \times n)$ 番目のサンプルに移動されることになる。このとき、ピッチ周波数を $p_1Hz$ とすると、h半音シフトした後のピッチ周波数を示すサンプル番号は第 $(p_1 \times 2^{h/12} \times N / f_s)$ 番目となる。

【0015】ここで、同じ人物が音程を変えて発音した声を分析したところ、音程が高くなるにつれピッチ周波数の倍音成分のレベルが比較的小さく、音程が低くなると倍音成分のレベルが大きくなり、豊富に出現することを見出した。そして、このピッチ周波数の倍音成分のレベルが再生される音声品質に影響を与えることが判ったので、周波数全体の移動後にこの倍音成分のレベルを操作して、高品質の音声にする。

【0016】ピッチ周波数抽出手段2において、抽出されたピッチ周波数が0である（ピッチ周波数が抽出されない）場合は（ステップ18→Yes）、倍音構造操作手段5に供給される音声信号は、何も操作せずにIFFT回路（逆フーリエ変換手段）6に出力される（ステップ22）。

【0017】ピッチ周波数抽出手段2において、抽出されたピッチ周波数が0でない（ピッチ周波数が存在する）場合は（ステップ18→No）、倍音構造操作手段5に供給される音声信号は、ピッチ周波数の倍音成分（ピッチ周波数の整数倍の周波数を示すサンプル）のレベルを操作する。即ち、周波数全体を高い方へシフト（シフト量 $\geq 1$ ）した場合には（ステップ19→Yes）、ピッチシフトした後の信号の倍音成分のレベルを減少させ（ステップ20）、周波数全体を低い方へシフト（シフト量 $< 1$ ）した場合には（ステップ19→No）、ピッチシフトした後の信号の倍音成分のレベルを増加させる（ステップ21）。本実施例では、共に10dBだけレベルを変化させることにしている。

【0018】例えば抽出されたピッチ周波数が200Hzであるとき、周波数全体を高い方へ3半音シフトした（ピッチシフト量が1倍以上）場合には、シフトした後のピッチ周波数は $200 \times 1.19$ Hzとなるので、シフトした後の音声信号の倍音成分は、 $200 \times 1.19 \times m$ （ $m$ は2以上の整数）Hzとなる。そして、この周波数を示すサンプル番号の実部及び虚部を各々 $10^{-0.5}$ 乗算して、約-10dBのレベル操作を行う。これを一般化すると、ピッチ周波数 $p_1$ Hzのときの $h$ 半音シフトした後の $m$ 倍音成分を示すサンプル番号は、第 $(m \times p_1 \times 2^{h/12} \times N / f_s)$ 番目となるので、このサンプル番号のデータの実部及び虚部を各々 $10^{-0.5}$ または $10^{0.5}$ を乗算することにより、 $\pm 10$ dBのレベル操作が可能となる。

【0019】この後、IFFT回路6に供給されて、逆フーリエ変換され、周波数領域から時間領域へ変換される（ステップ22）。IFFT回路6により時間領域の信号に変換された音声信号は、フィルタ7に供給されて再び第0番～第999番のサンプルは正弦波の窓関数で時間窓で切り出され、第3096番～第4095番のサンプルは余弦波の窓関数で時間窓で切り出され、その他のサンプルは1の窓関数でフィルタをかけられて出力される（ステップ23）。そして、最初の音声信号の第3096番～第4095番のサンプルデータを図示せぬメモリ等に格納しておき、第0番～第3095番のサンプルデータをD/A変換器（図示せぬ）などへ出力する。

【0020】次に入力される音声信号のデータは、最初の音声信号の第3096番のサンプルから4096サンプル分を読み出して、上記と同様の処理を行う。そして、図3に示すように、フィルタ7から出力される音声信号に対して先に格納していた最初の音声信号の第30

96番～第4095番のサンプルデータを加算する（ステップ24）と共に、このサンプルデータの最後の部分1000サンプルのデータを図示せぬメモリ等に格納する（ステップ25）。この様に、正弦波または余弦波の窓関数で時間窓で切り出される前後1000サンプル分のデータが重なるように切り出して、重なる部分のデータを加算しながら出力していく（ステップ26）。そして、フレーム番号 $i$ に1を加算し（ステップ27）、入力される音声信号がなくなるまで、これらの処理を繰り返す。

【0021】なお、上記実施例での処理区間は4096サンプルとしているが、これ以外のサンプル数でも良いのは勿論である。しかしながら、種々の実験を行った結果、1サンプル当たり10Hz～25Hz程度となるように処理区間を設定するのが音質上最も良いことが判った。そして、フーリエ変換等のデジタル処理を行うことを考慮すると、処理区間は2の $n$ 乗サンプルにするのが良い。したがって、上記実施例のようにサンプリング周波数44.1kHzの音声データの場合は、2048サンプル（21.5Hz/1サンプル）または4096サンプル（10.8Hz/1サンプル）とするのが良く、MPEG2オーディオ等で使用されるサンプリング周波数22.05kHzの音声データの場合は、1024サンプル（21.5Hz/1サンプル）または2048サンプル（10.8Hz/1サンプル）とするのが良い。

【0022】実際に、サンプリング周波数44.1kHzの音声データについて、処理区間を512、1024、2048、4096、8192の各サンプルで実験したところ、512サンプルでは音程が一つに定まらず、1024サンプルでは音質が非常に悪かった。そして、8192サンプルでは所望の音程にはなったものの、ディレイがかかったような2重の音声となってしまう、処理区間は2048または4096サンプルのときが最も高音質の結果を得ることができた。

#### 【0023】

【発明の効果】本発明の音程変換装置は、音声信号のピッチ周波数を抽出して、フーリエ変換した後に全周波数帯域を高域側または低域側にシフトした音声信号のピッチ周波数の倍音の構造を操作してから逆フーリエ変換することにより、周波数領域で倍音成分の特徴を維持したまま全周波数帯域をシフトしているので、従来に比べ簡単な回路構成で処理時間も比較的短く、しかも音質の劣化がなくて個人の声の特徴を維持したままの自然で高品質な音声音程変換が可能となるという効果がある。

#### 【図面の簡単な説明】

【図1】本発明の音程変換装置の一実施例を示すブロック図である。

【図2】本発明の音程変換装置の一実施例を示すフローチャート図である。

【図3】本発明の音程変換装置の一実施例の時間窓での

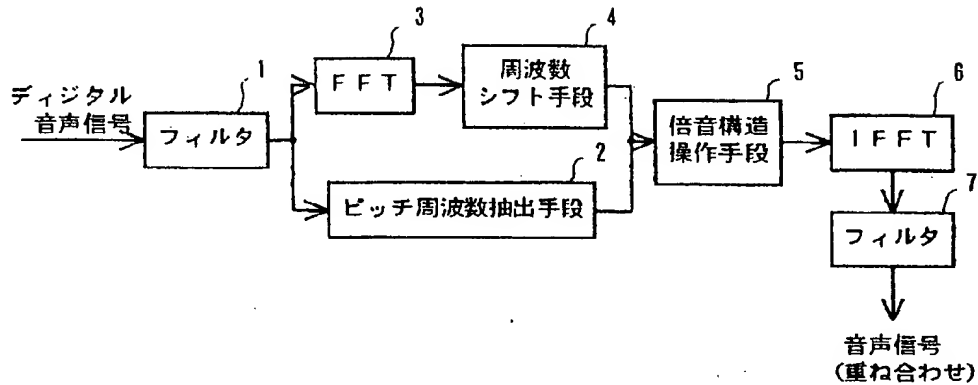
切り出しと重ね合わせを説明するための図である。

【符号の説明】

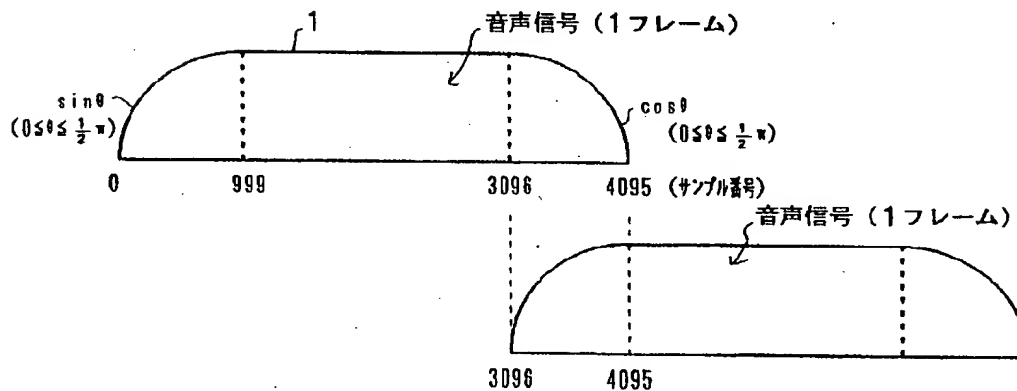
- 1 フィルタ (分割手段)
- 2 ピッチ周波数抽出手段
- 3 F F T 回路 (フーリエ変換手段)

- 4 周波数シフト手段
- 5 倍音構造操作手段
- 6 I F F T 回路 (逆フーリエ変換手段)
- 7 フィルタ

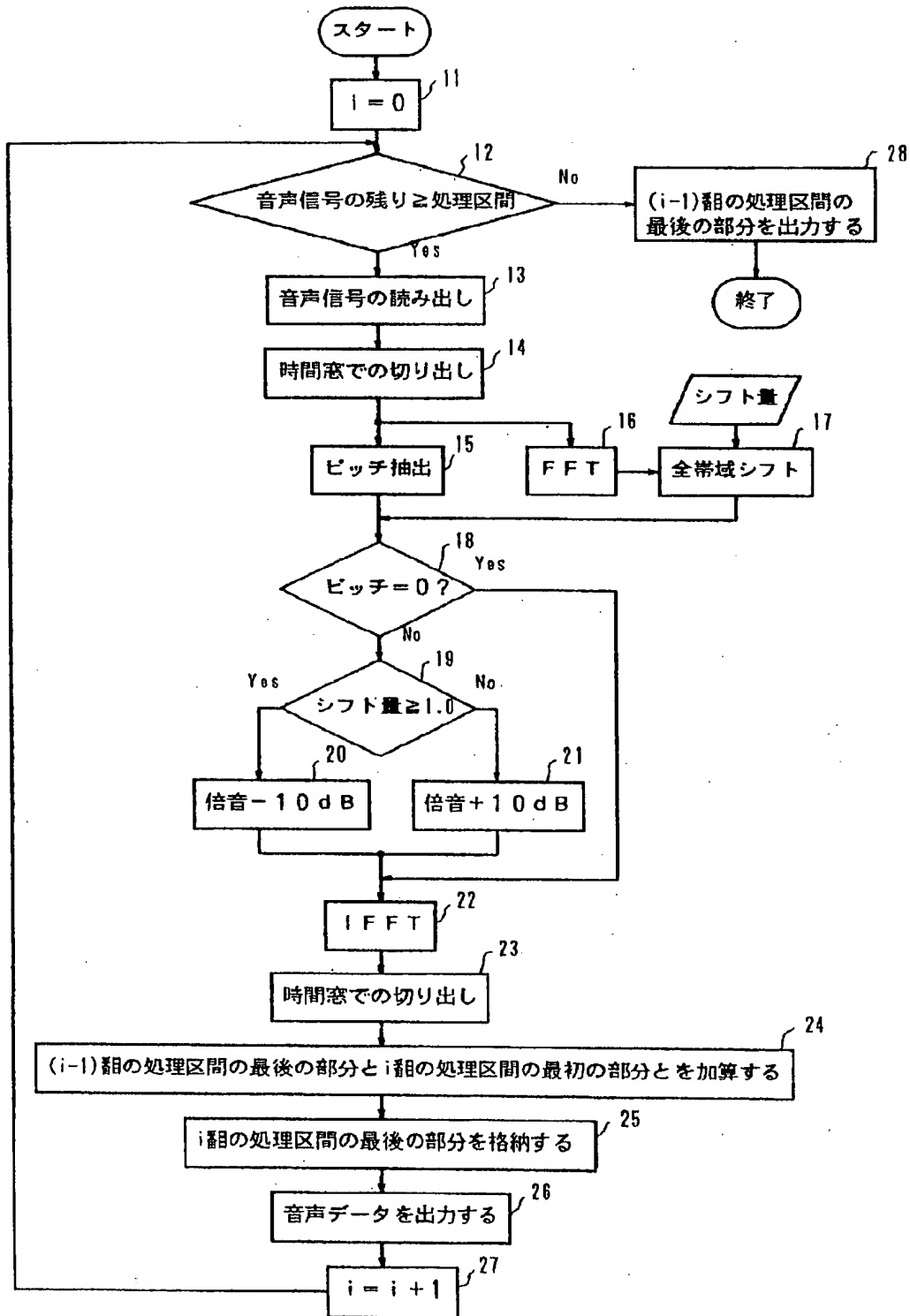
【図 1】



【図 3】



【図 2】





## 【手続補正書】

【提出日】平成8年9月4日

## 【手続補正1】

【補正対象書類名】明細書

【補正対象項目名】請求項2

【補正方法】変更

## 【補正内容】

【請求項2】前記分割手段は、デジタル入力された音声信号を所定時間のフレームに切り出すと共に、このフレームの最初から10～35msecまでのデータを正弦波の1/4周期分の時間窓で切り出し、このフレームの最後から10～35msecまでのデータを余弦波の1/4周期分の時間窓で切り出すことを特徴とする請求項1記載の音程変換装置。

## 【手続補正2】

【補正対象書類名】明細書

【補正対象項目名】0012

【補正方法】変更

## 【補正内容】

【0012】そして、このフィルタ1における正弦波および余弦波による時間窓での切り出しは、200～2000サンプル幅の任意サンプル幅の区間で種々実験したところ、音源によって多少の変化はあるが、ほとんどの音源で500～1500サンプル（約10～35msec）幅の間が最適な区間になることが判ったので、この実施例では1000サンプル（約23msec）幅で正弦波および余弦波による時間窓での切り出しを行っている。なお、この切り出し区間のサンプル数（500～1

500サンプル）は、フレームサンプル数の半分以下の範囲で変更可能である。このフィルタ1により切り出された音声信号は、ピッチ周波数抽出手段2に供給されて、自己相関関数やケプストラム法等によりピッチ周波数（ピーク周波数のうち最も低い周波数（基本周波数）を示すサンプル）が抽出される（ステップ15）。また、フィルタ1より出力された音声信号は、FFT回路（フーリエ変換手段）3にも供給されてフーリエ変換を施され、時間領域の信号から周波数領域の信号へ変換される（ステップ16）。

## 【手続補正3】

【補正対象書類名】明細書

【補正対象項目名】0013

【補正方法】変更

## 【補正内容】

【0013】このとき、時間領域に対応していた各サンプルは、各周波数に対応し、サンプル番号と周波数とが対応することになる。即ち、サンプリング周波数 $f_s$ の音声信号データをN個のサンプル毎に切り出して処理する場合、FFT回路3から出力される信号の周波数 $pHz$ を示すサンプル番号は第 $(p \times N / f_s)$ 番目となる。本実施例の場合、サンプリング周波数44.1kHzの音声信号データに対して4096サンプル毎に切り出しているため周波数 $pHz$ を示すサンプル番号は第 $(p \times 4096 / 44100)$ 番目となる（小数点以下四捨五入）。

【公報種別】特許法第17条の2の規定による補正の掲載

【部門区分】第6部門第2区分

【発行日】平成13年2月9日(2001.2.9)

【公開番号】特開平9-185392

【公開日】平成9年7月15日(1997.7.15)

【年通号数】公開特許公報9-1854

【出願番号】特願平7-353508

【国際特許分類第7版】

G10L 21/04

G10H 1/00

G10K 15/04 302

【FI】

G10L 3/02 A

G10H 1/00 B

G10K 15/04 302 D

【手続補正書】

【提出日】平成12年3月24日(2000.3.24)

【手続補正1】

【補正対象書類名】明細書

【補正対象項目名】請求項2

【補正方法】変更

【補正内容】

【請求項2】前記分割手段は、デジタル入力された音

声信号を所定時間のフレームに切り出すと共に、このフレームの最初の0～35msec(上限は10～35msecの間で変更可能)のデータを正弦波の1/4周期分の時間窓で切り出し、このフレームの最後の0～35msec(上限は10～35msecの間で変更可能)のデータを余弦波の1/4周期分の時間窓で切り出すことを特徴とする請求項1記載の音程変換装置。